

A new simulation approach for HPC interconnects

Ronald Luijten (lui@zurich.ibm.com)
IBM Rueschlikon
11 March 2009





Andreas Doering*

Krzysztof Kryszczuk*

Phillip Stanley Marbell

Gergely Paljak

Alessandra Scicchitano

German Rodriguez Herrera

HansPeter Ineichen (borrowed from scs)

Rolf Clauberg

Wolfgang Denzel

Cyriel Minkenberg

Mitch Gusat

Francois Abel *

Patricia Sagmeister*

Maria Gabrani*

Lydia (Yiyu) Chen

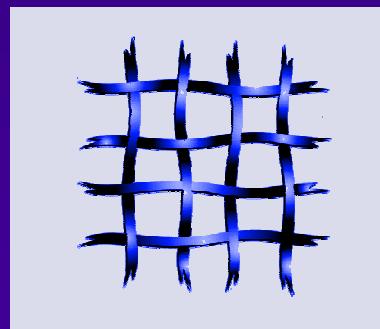
Florian Auernhammer

Alternative title:

A Journey from Telco Switching to Computer Interconnects...



Ended 2003

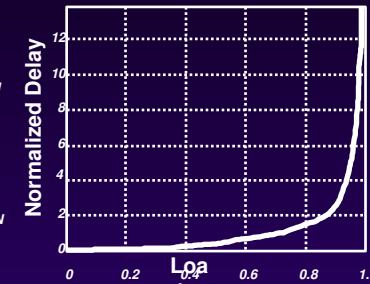
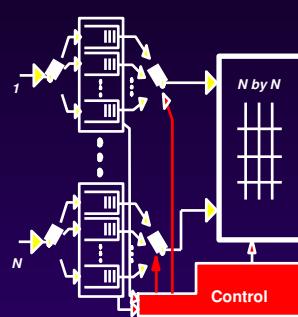


2009

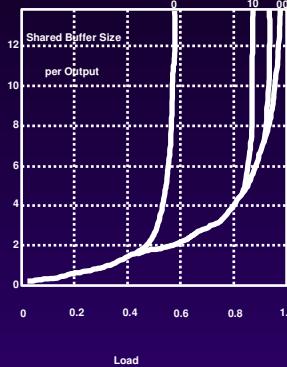
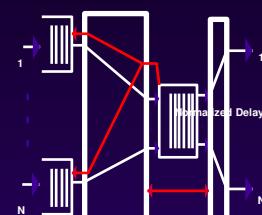
PRIZMA's Robust Performance



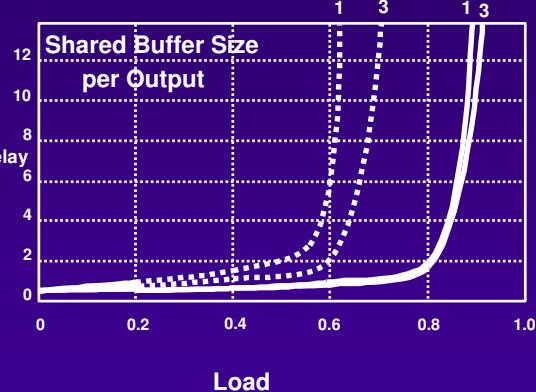
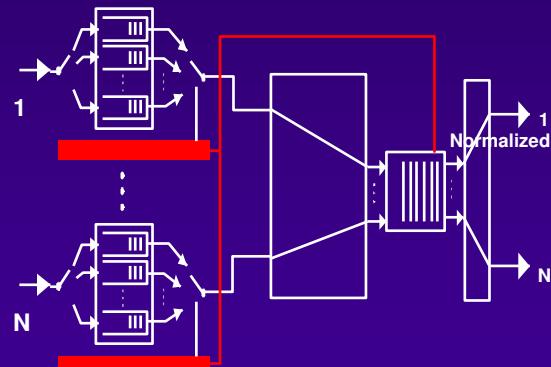
Virtual Output Queuing



Shared output Queueing



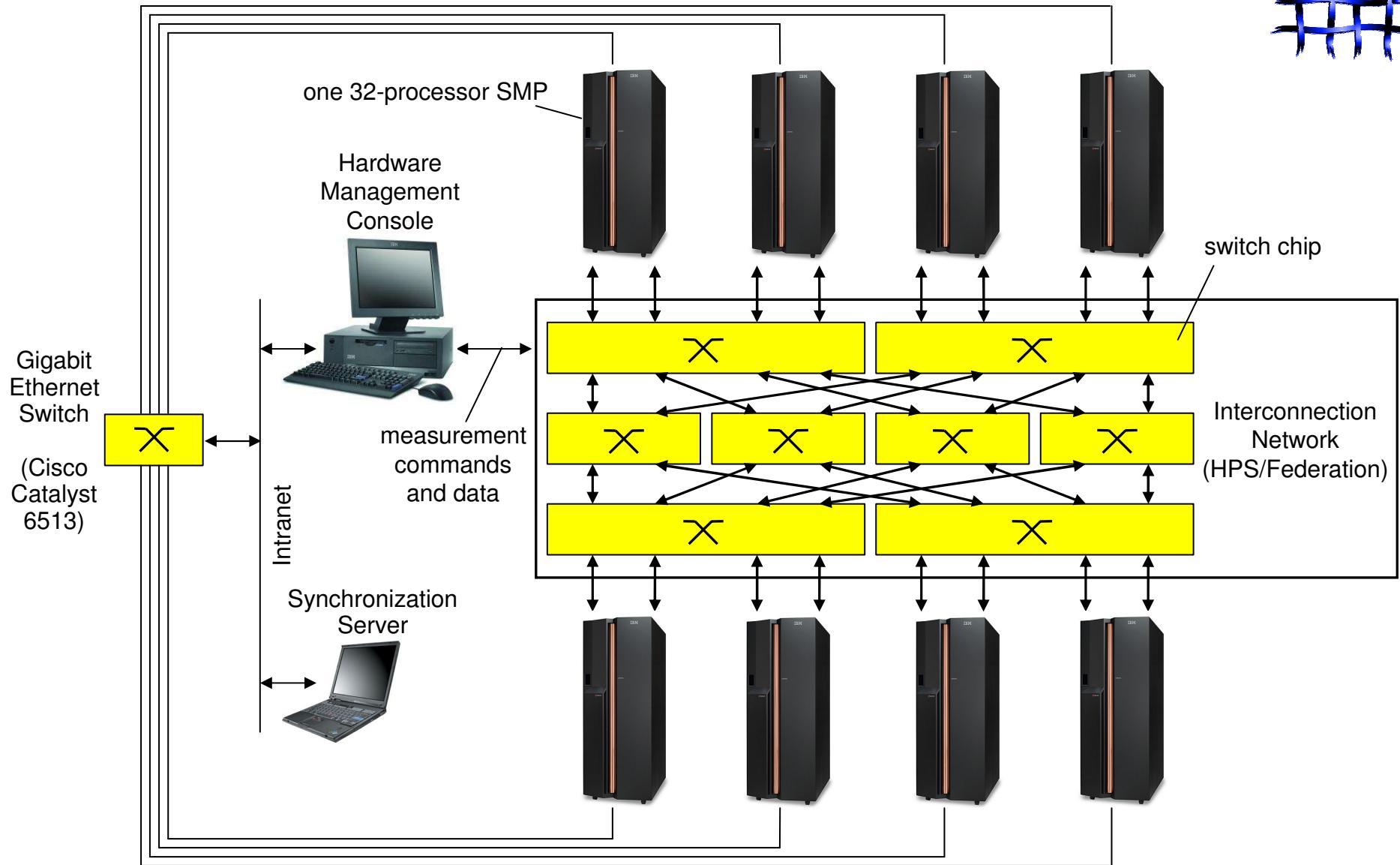
Virtual Output Queuing with Output Buffering

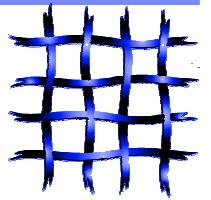


IP

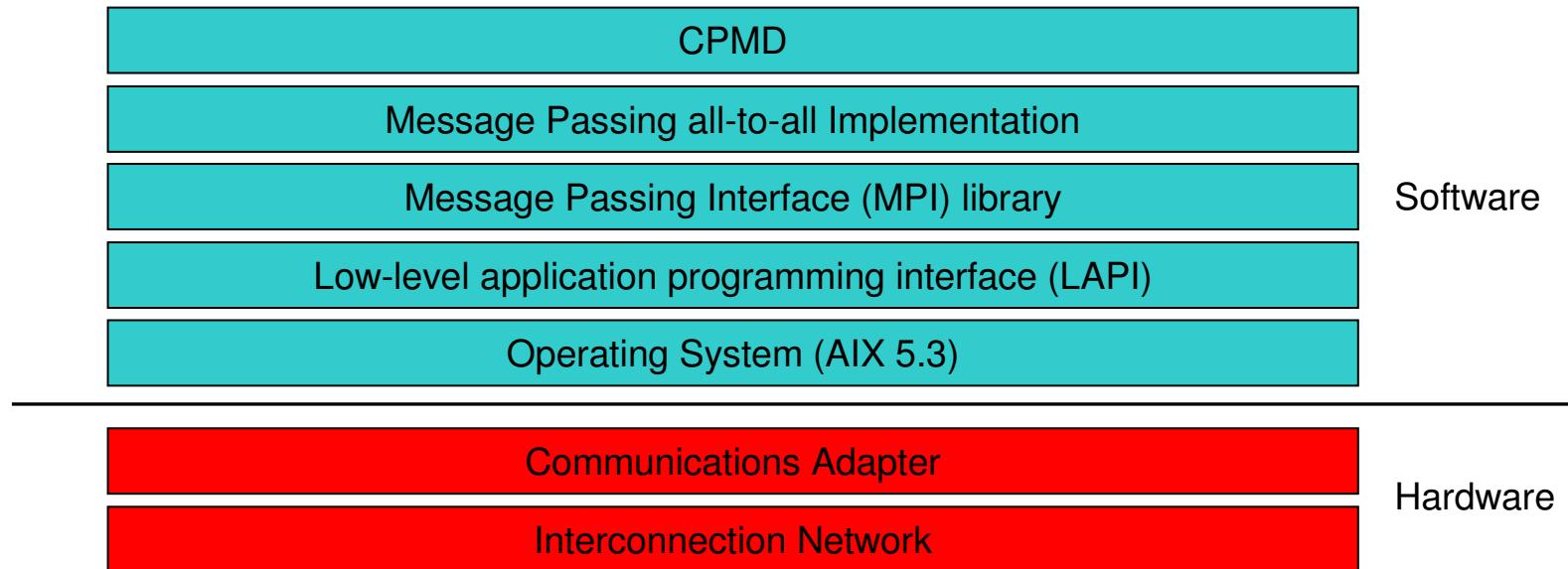
- Robust switch performance for extreme traffic conditions
- memory used for contention resolution only - not for performance
- more scalable solution compared to VOQ - crossbar approach
- Documented: IEEE Comm.Mag. Dec. 2000 - Best Comm. Mag. Paper 2000

Measurement on ZRL Regatta w/ Federation

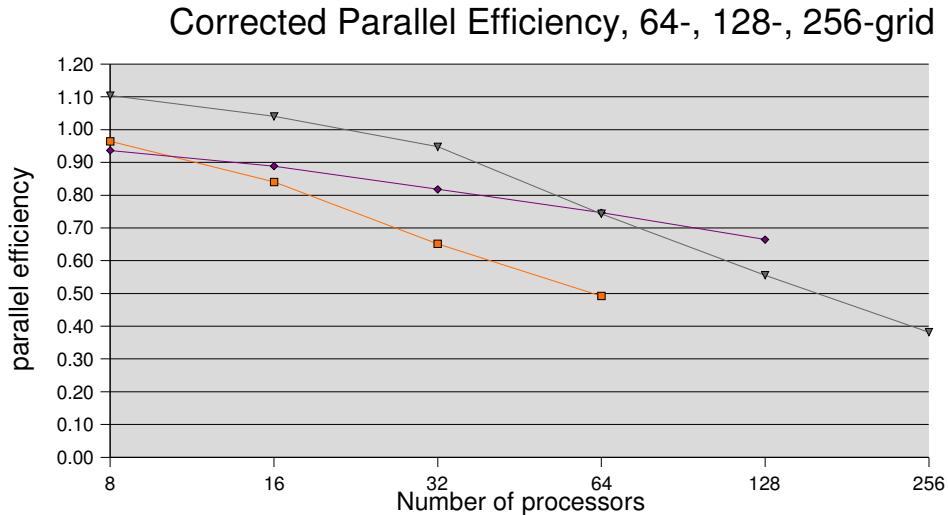
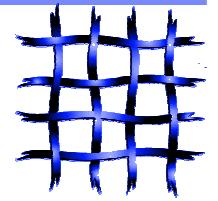




Regatta CPMD Communication stack

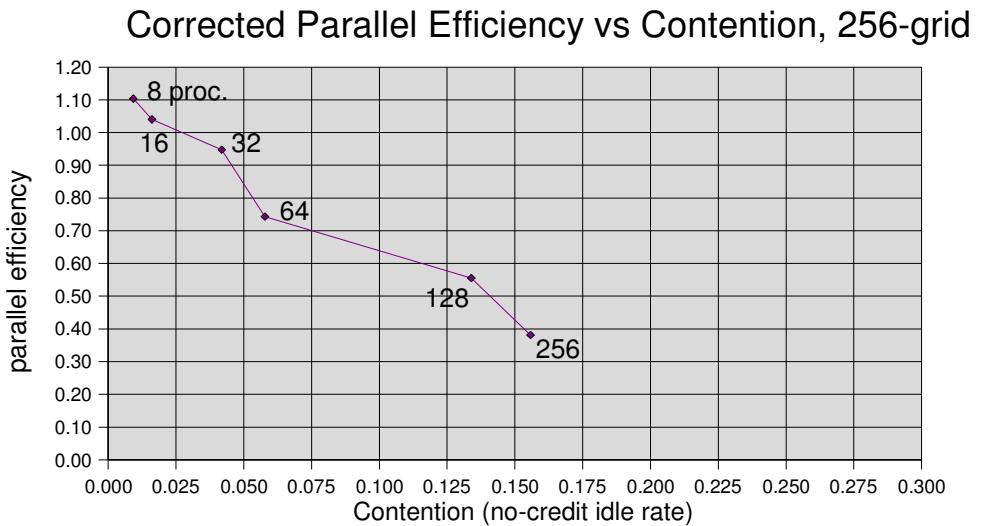


Measurement Results (Parallel Efficiency)

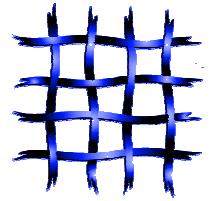


- Measurement config:
 - IBM MPI
 - LAPI
 - IBM SMA3 adapter
 - user-space mode
 - HPS
- Ideally: parallel efficiency is constant (1.00)

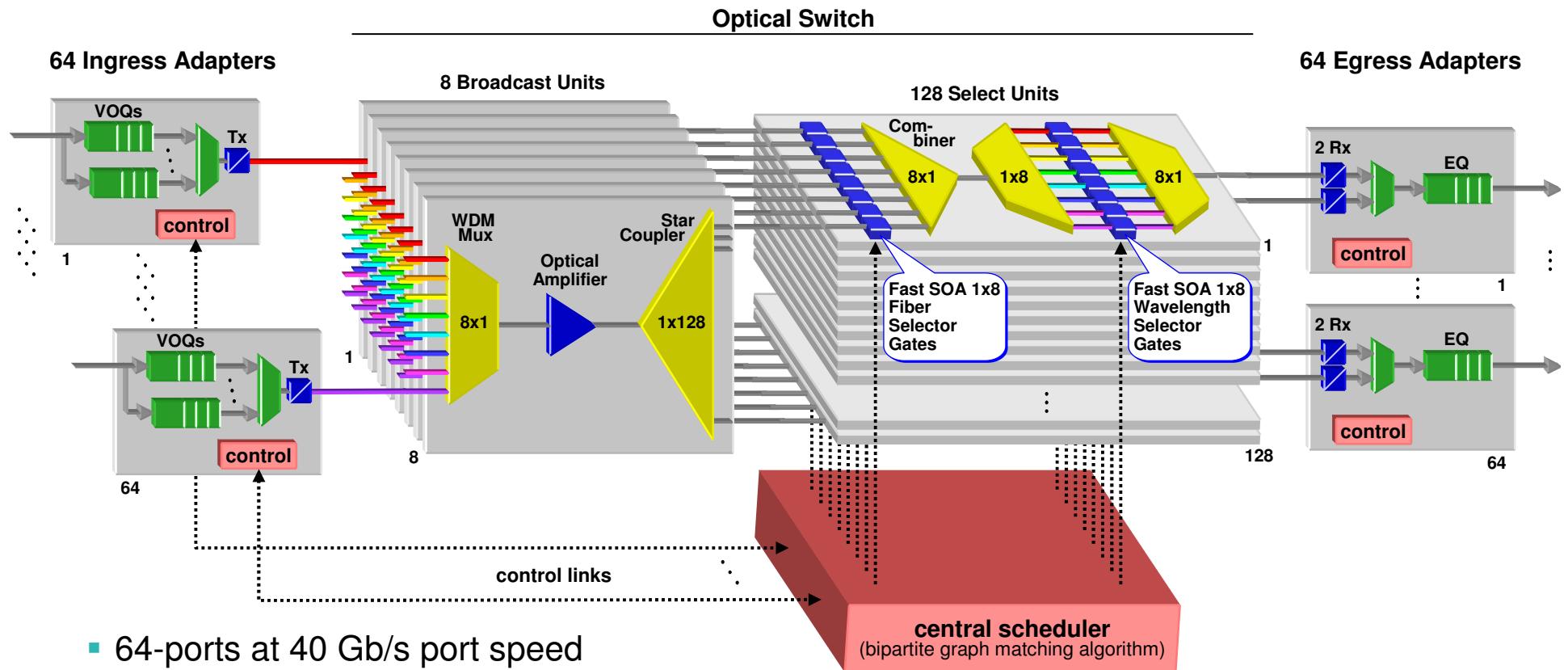
- Strong evidence that **contention in the interconnect causes decrease in parallel efficiency**
- Interconnect only ~20% used
- **Suspect effect of interconnect contention on inefficiency in the software layers**



OSMOSIS Demonstrator System - Overview

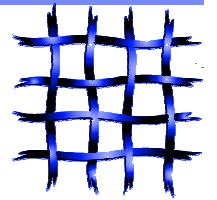


(Optical Shared MemOry Supercomputer Interconnect System)

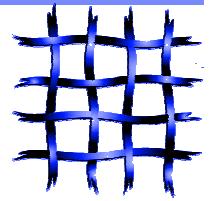


- 64-ports at 40 Gb/s port speed
- Broadcast-and-select architecture
- Combination of wavelength and space division multiplexing
- fast switching based on SOAs

Insights gained 3 years after our transition



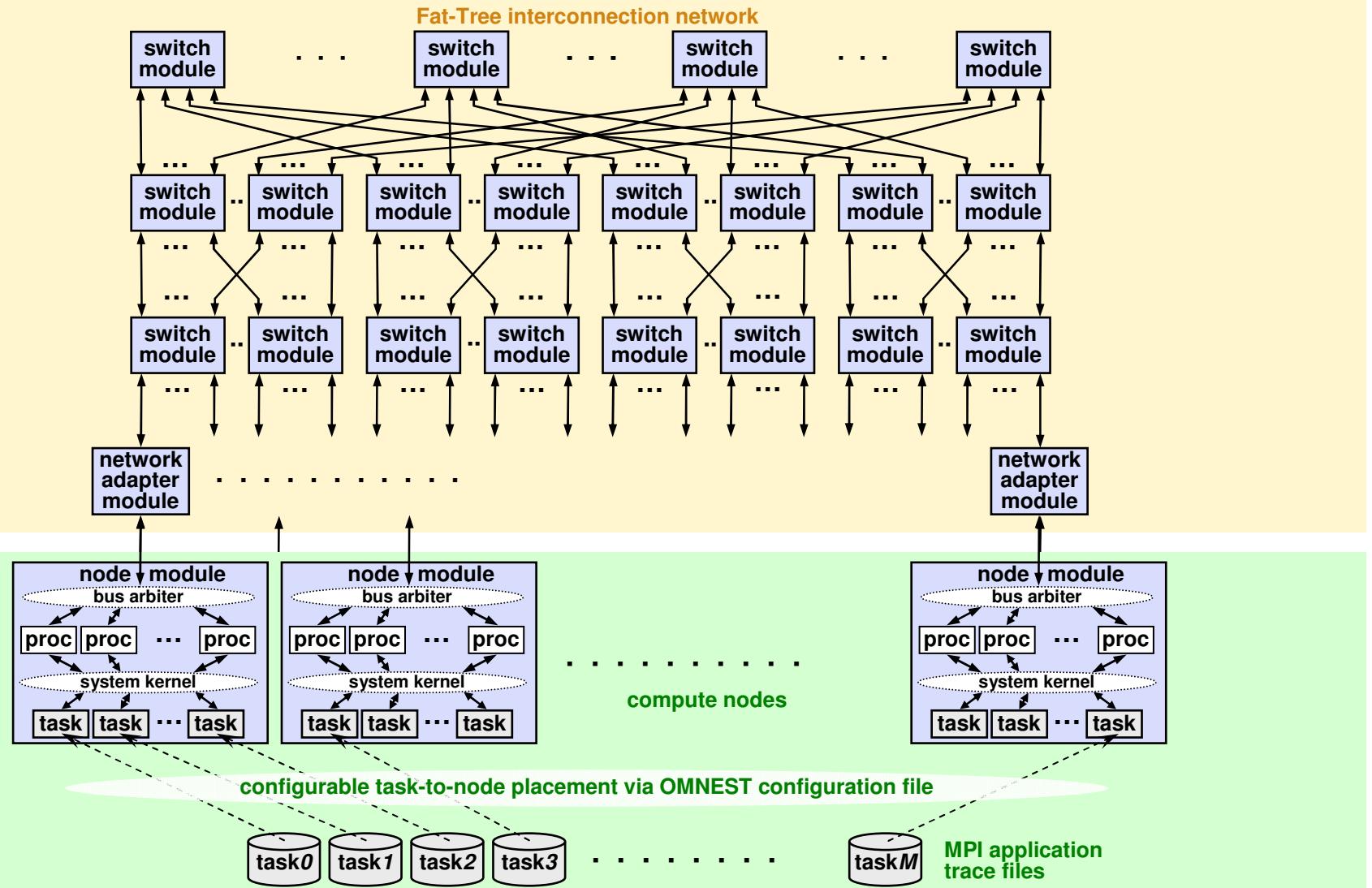
PERCS



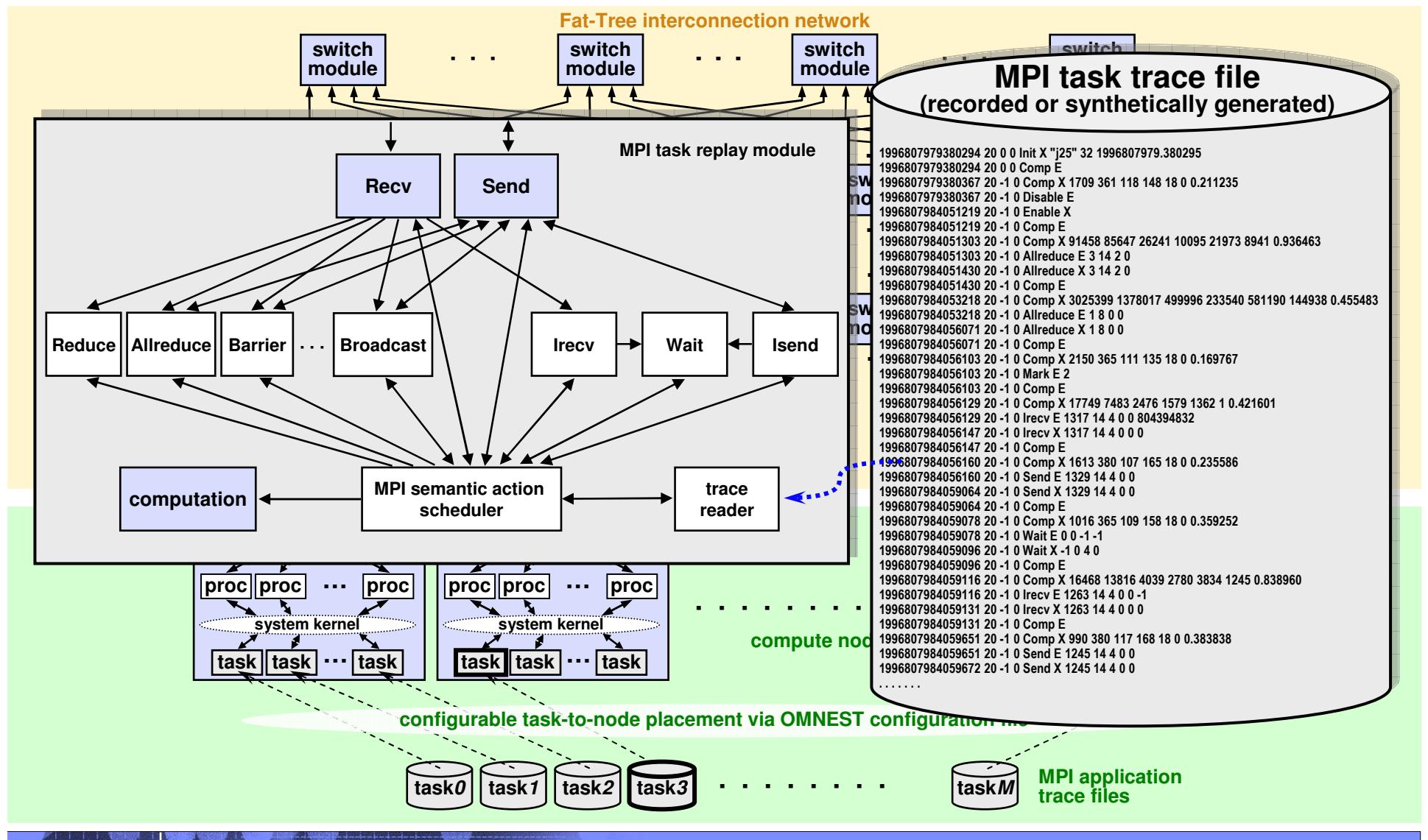
HPCs
PERCS

A 3D-style blue arrow points diagonally upwards from the bottom left towards the top right, positioned behind the text "HPCs PERCS".

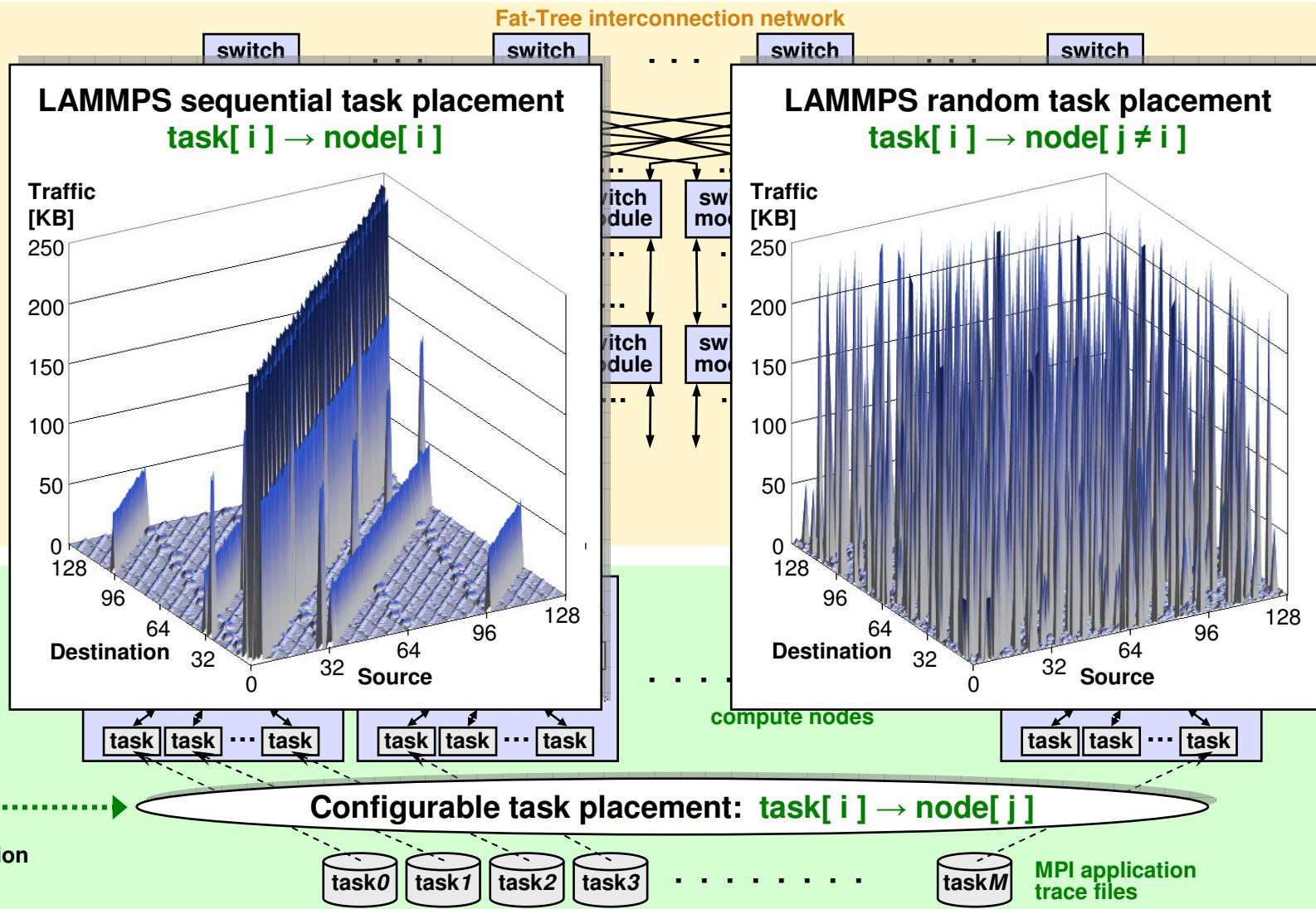
PERCS phase-2 full-system model

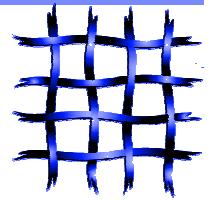


MARS trace replay



MARS task placement

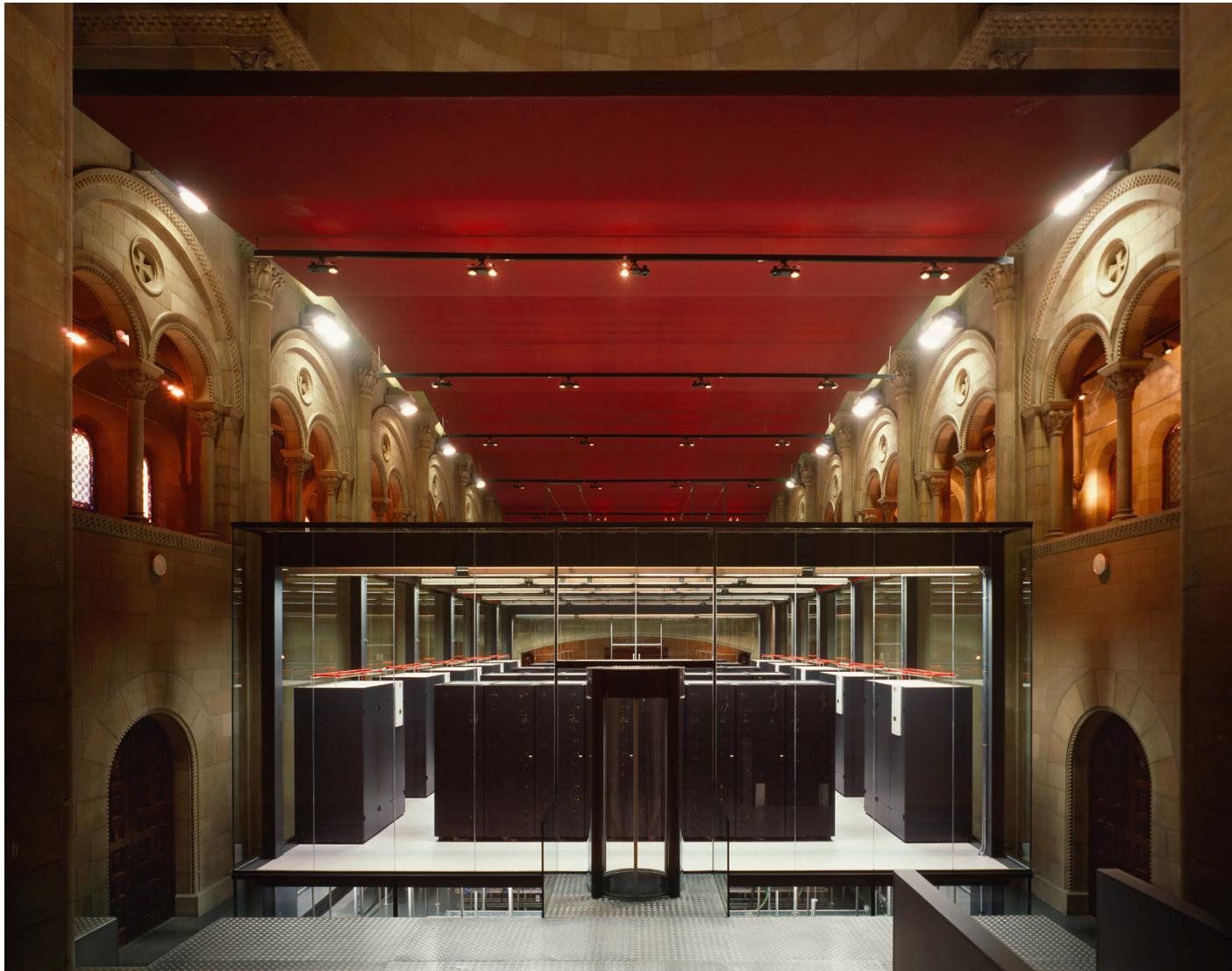
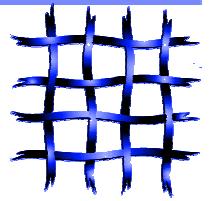


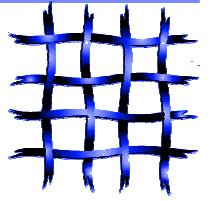


Good new direction... However

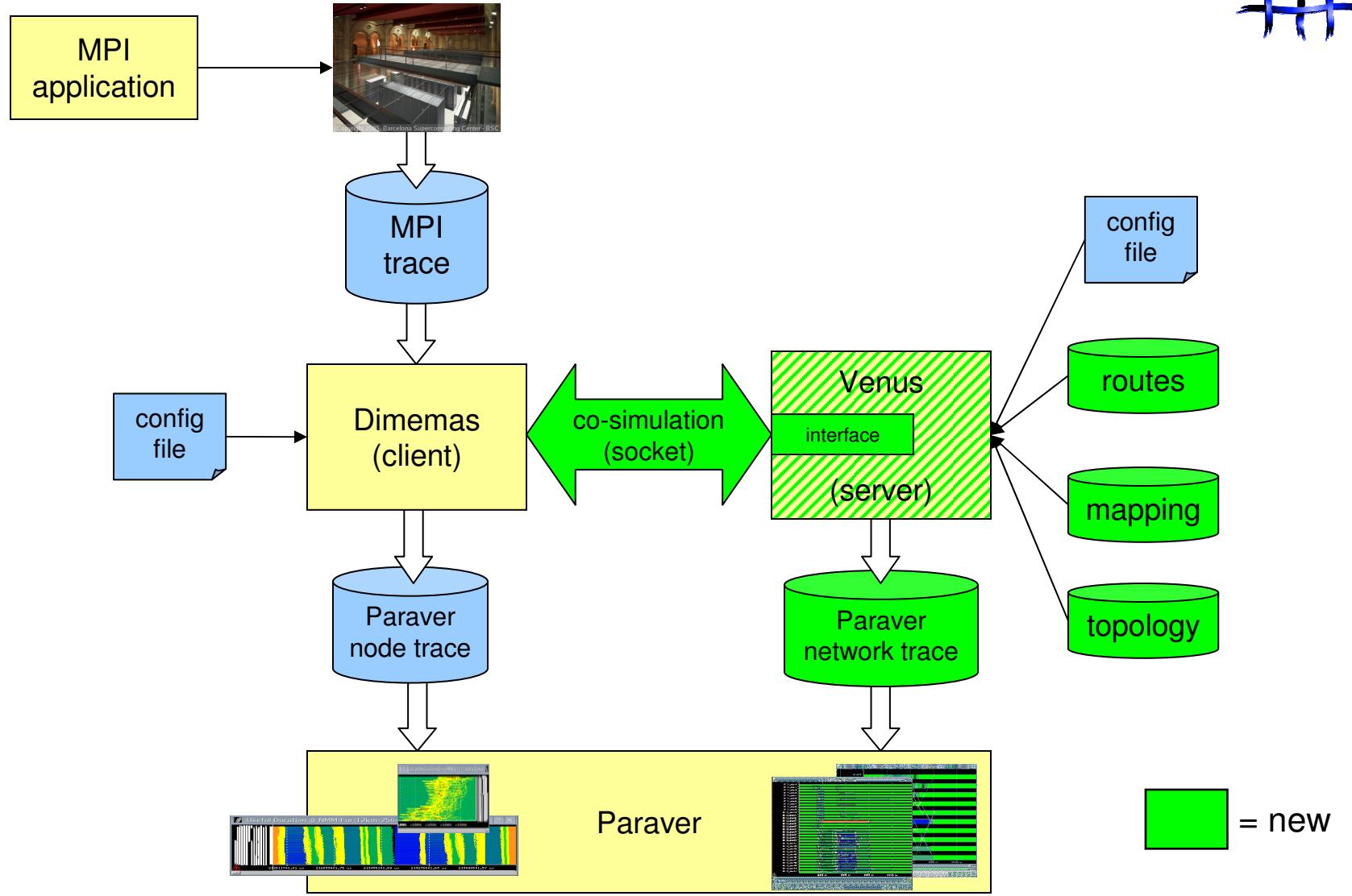
Where do I get good traces?

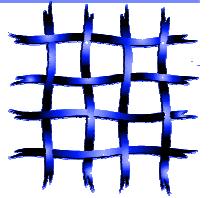
Mare Incognito



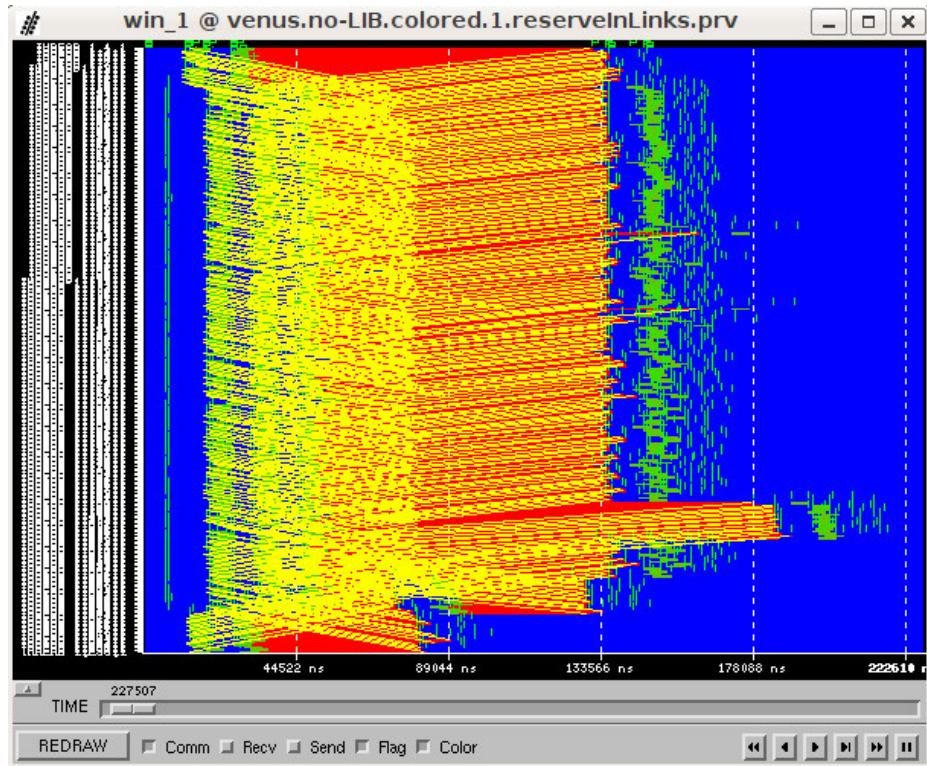


Tool Enhancement & Integration



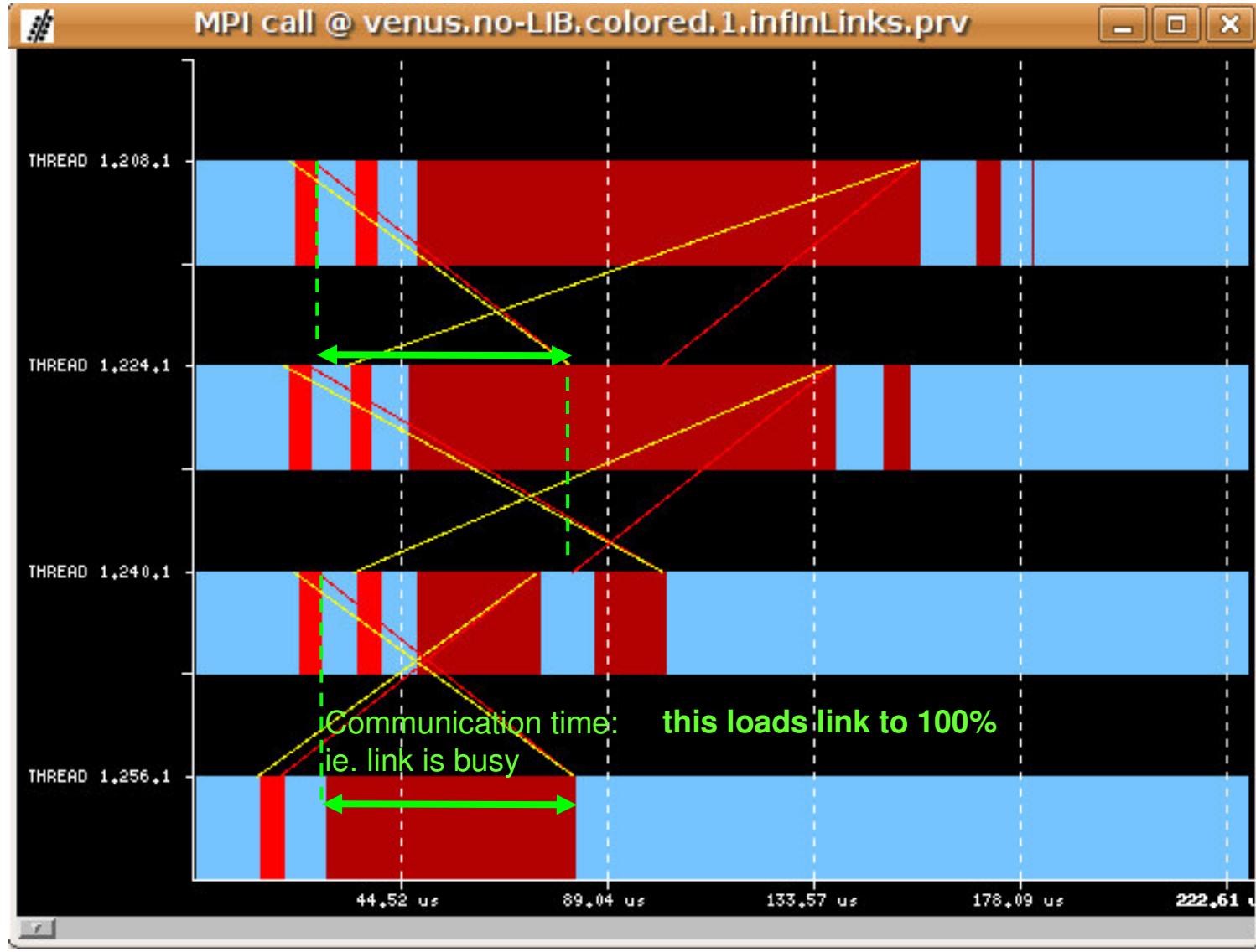


WRF application example

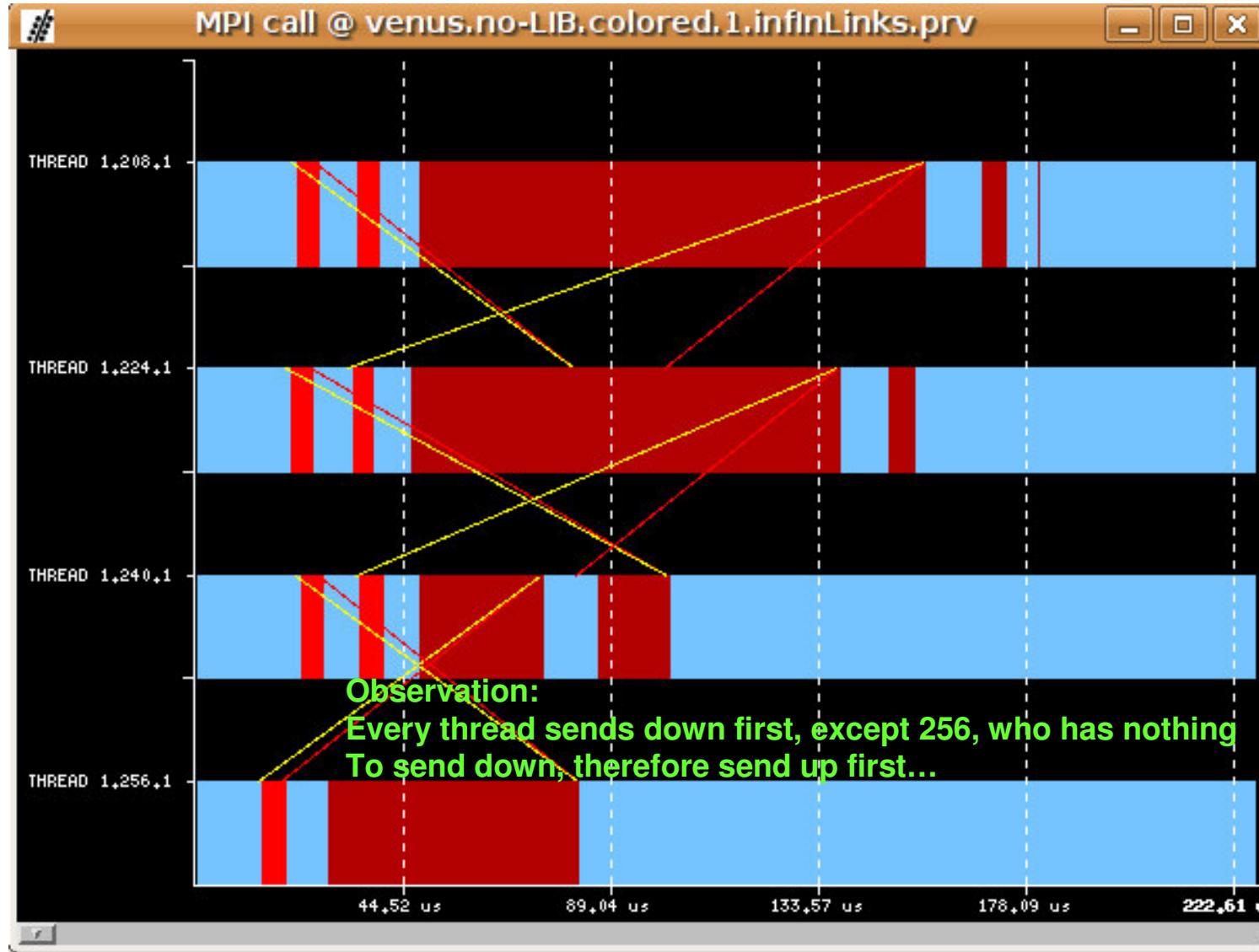
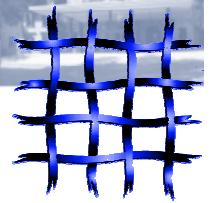


- WRF: 256 threads, sim time ~ 2hours
 - Routes with no conflicts
 - Simulation by Venus
 - Load Balanced
- Some threads take longer:
 - Dependency chains

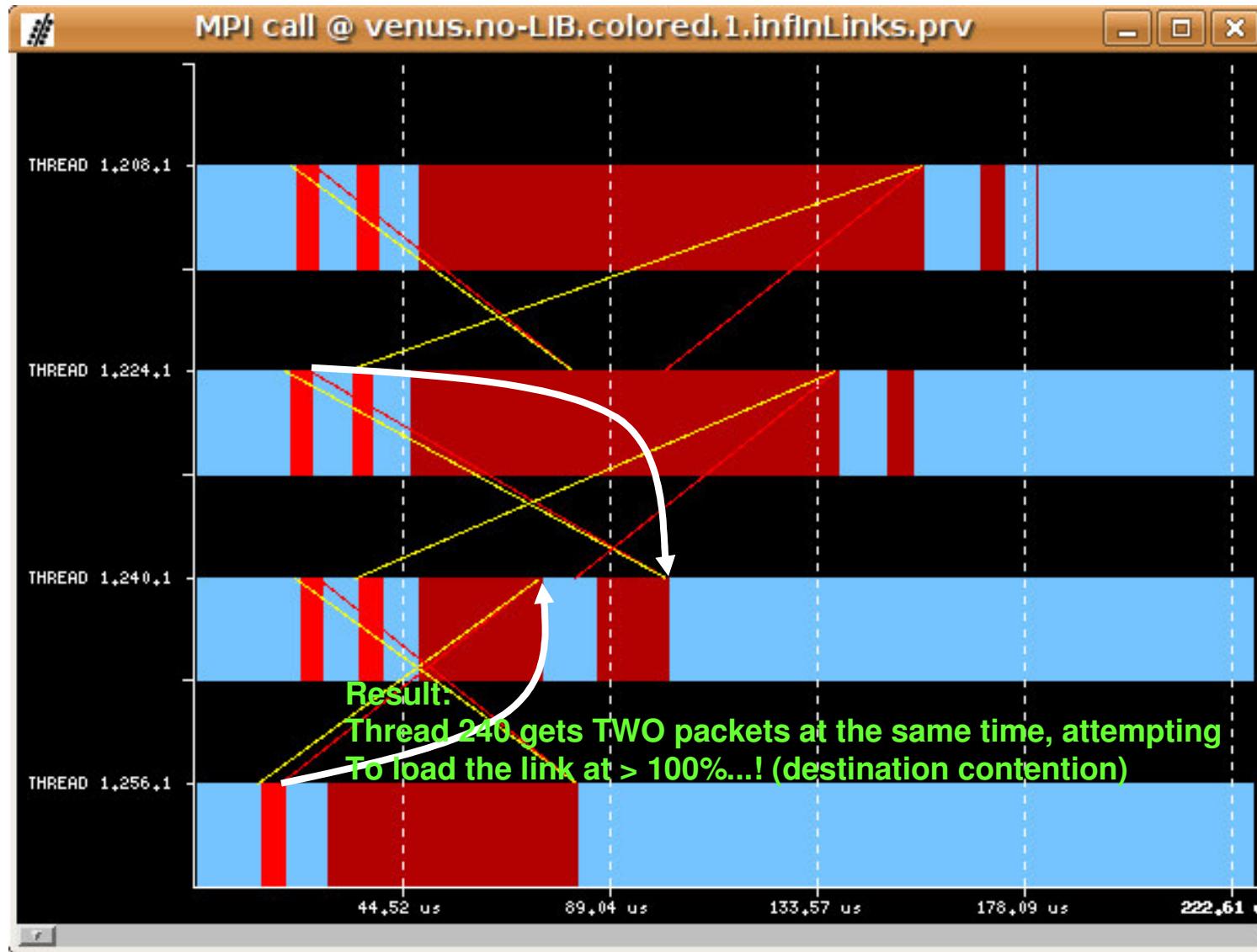
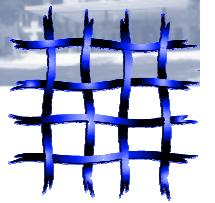
Communication Dependency chains



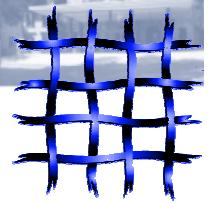
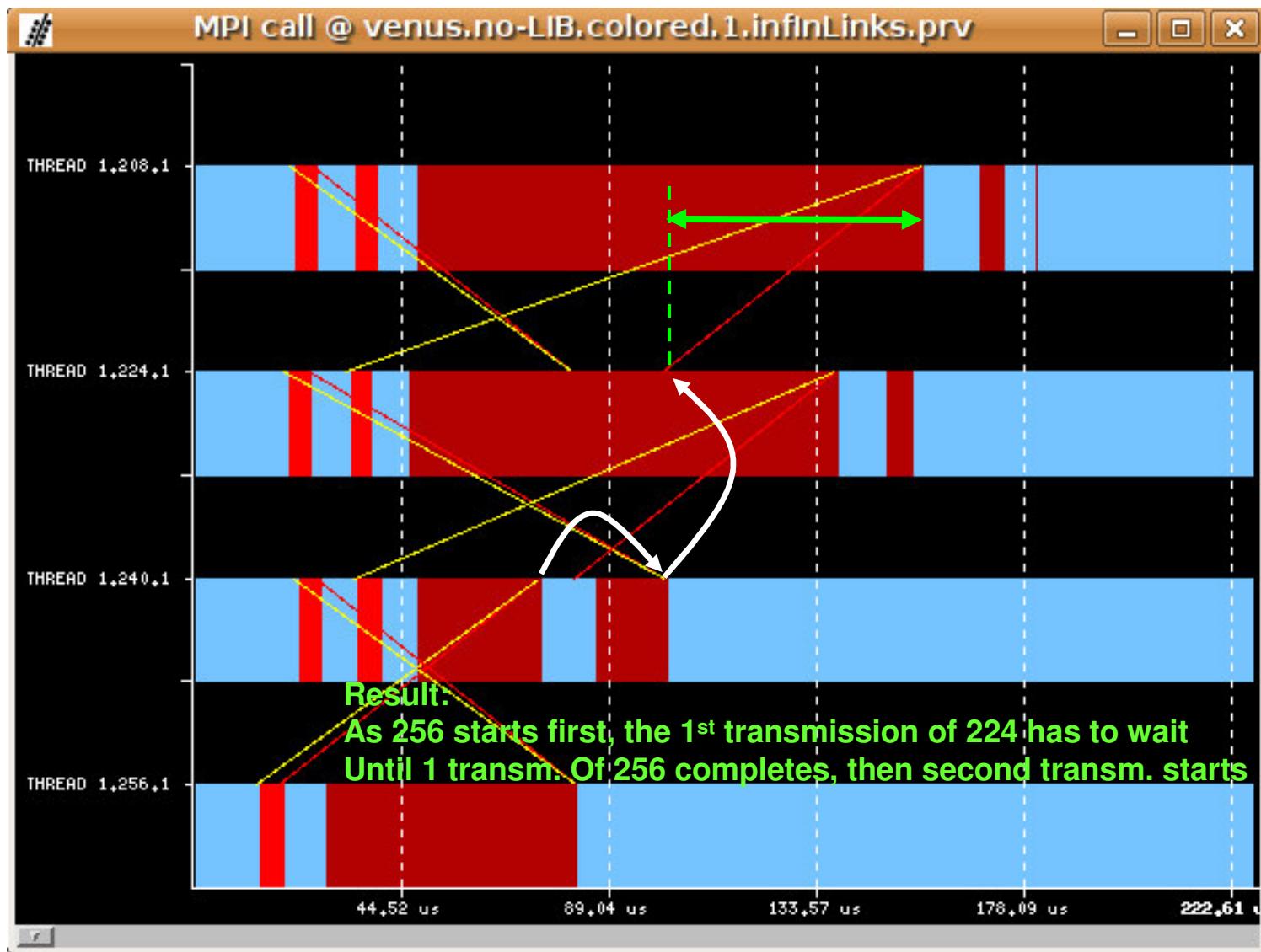
Communication Dependency chains



Communication Dependency chains



Communication Dependency chains



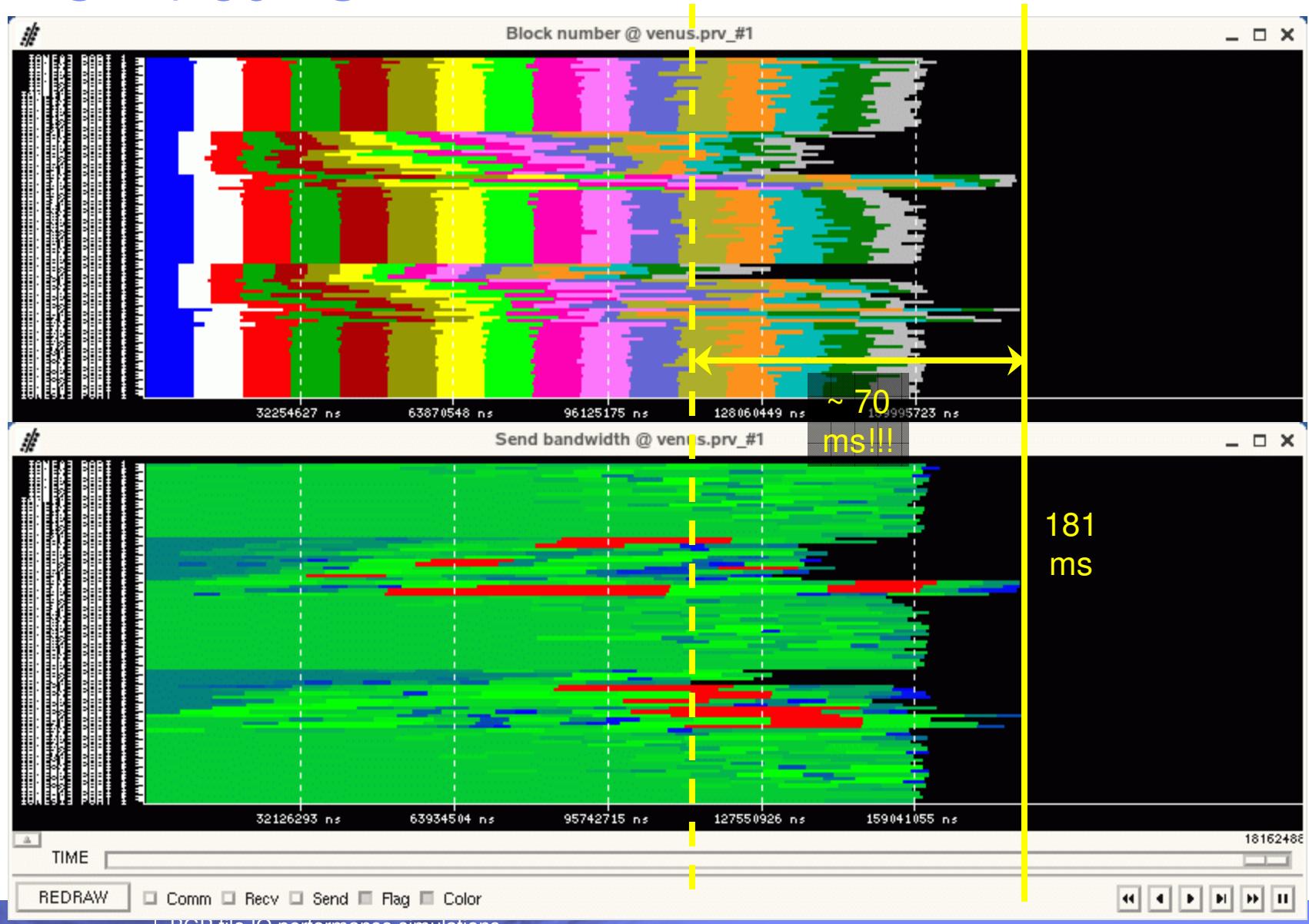
We also study file I/O performance...

Checkpointing on a BG/P machine, using GPFS

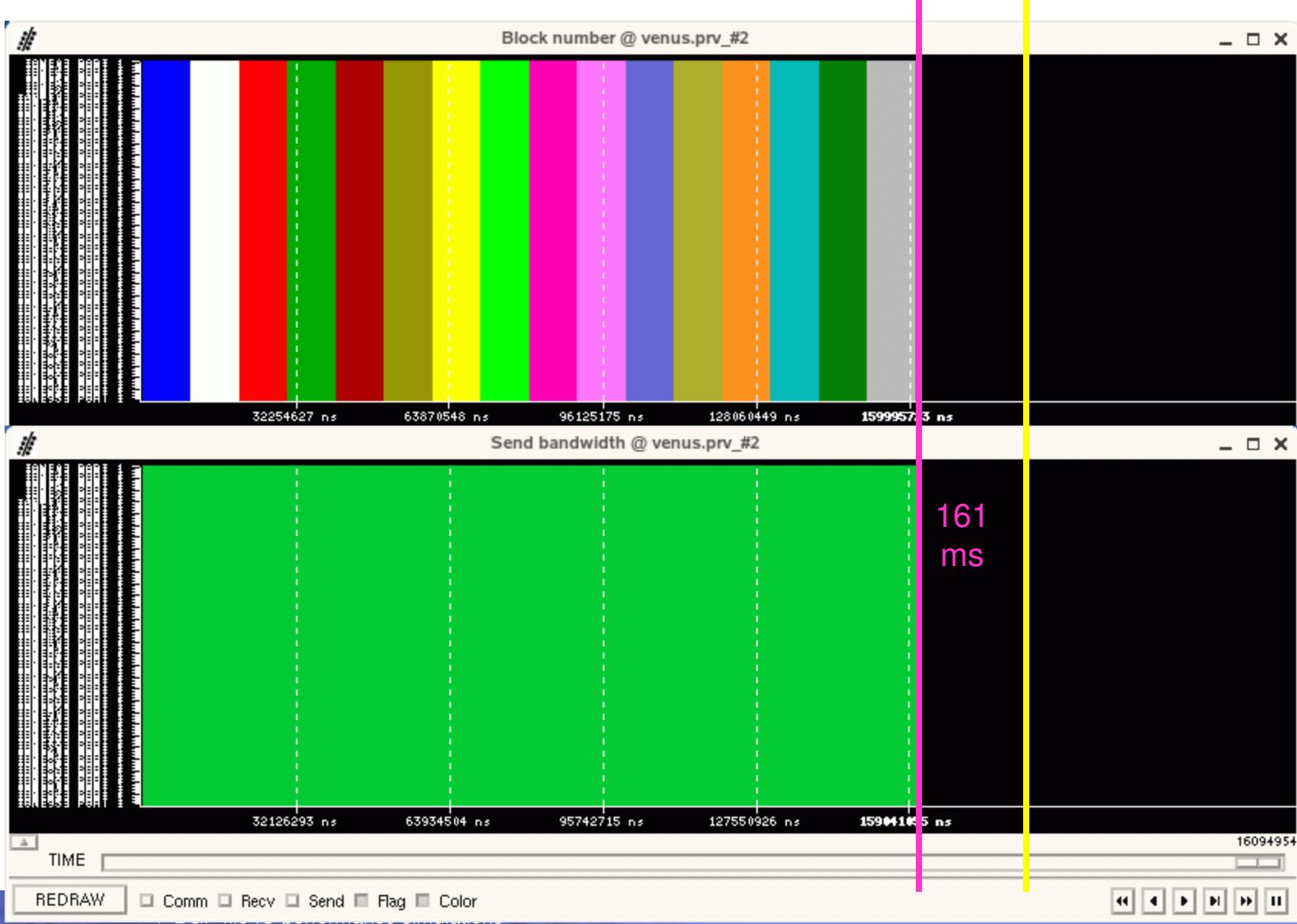
Simulation parameters

- Network
 - 128 nodes total
 - Three-level fat tree
 - 8-port switches: $32+32+16 = 80$ switches
- Nodes
 - One compute node per IO node (ION)
 - 16 blocks of 4 MB each per node (64 MB per node)
 - ACK size = 512 B
- Myrinet
 - Flit size = 256 B
 - Flit duration = 204.8 ns (10 Gb/s)
 - Switch buffer size = 4 KB/port
 - No overhead (raw BW = user BW)
 - 4 KB segments; interleaved in groups of 4 segments by adapter
 - All latencies zero (link, switch, adapter)
- Three configurations
 - 92 ION + 36 FS
 - 96 ION + 32 FS
 - 100 ION + 28 FS

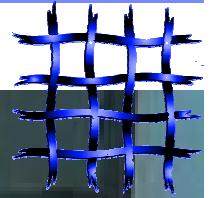
92 ION + 36 FS



96 ION + 32 FS



Insights gained 5 years after our transition



- Compute vs communicate
- What am I measuring, really?
- If a porsche were your hpc

Literature

- “**A Framework for End-to-end Simulation of High performance Computing Systems**”, Wolfgang E. Denzel, Jian Li, Peter Walker, Yuho Jin. Simutools 2008
- “**Trace-driven Co-simulation of High-Performance Computing Systems using OMNeT++**”, Cyriel Minkenberg, Germán Rodriguez Herrera. Simutools 2009